

## 1. Preparing our dataset

*These recommendations are so on point! How does this playlist know me so well?*



Over the past few years, streaming services with huge catalogs have become the primary means through which most people listen to their favorite music. But at the same time, the sheer amount of music on offer can mean users might be a bit overwhelmed when trying to look for newer music that suits their tastes.

For this reason, streaming services have looked into means of categorizing music to allow for personalized recommendations. One method involves direct analysis of the raw audio information in a given song, scoring the raw data on a variety of metrics. Today, we'll be examining data compiled by a research group known as The Echo Nest. Our goal is to look through this dataset and classify songs as being either 'Hip-Hop' or 'Rock' - all without listening to a single one ourselves. In doing so, we will learn how to clean our data, do some exploratory data visualization, and use feature reduction towards the goal of feeding our data through some simple machine learning algorithms, such as decision trees and logistic regression.

To begin with, let's load the metadata about our tracks alongside the track metrics compiled by The Echo Nest. A song is about more than its title, artist, and number of listens. We have another dataset that has musical features of each track such as `danceability` and `acousticness` on a scale from -1 to 1. These exist in two different files, which are in different formats - CSV and JSON. While CSV is a popular file format for denoting tabular data, JSON is another common file format in which databases often return the results of a given query.

Let's start by creating two pandas `DataFrames` out of these files that we can merge so we have features and labels (often also referred to as `X` and `y`) for the classification later on.

In [1]:

```
import pandas as pd

# Read in track metadata with genre labels
tracks = pd.read_csv('fma-rock-vs-hiphop.csv')

# Read in track metrics with the features
echonest_metrics = pd.read_json('echonest-metrics.json', precise_float = True)

# Merge the relevant columns of tracks and echonest_metrics
echo_tracks = echonest_metrics.merge(tracks[['track_id', 'genre_top']], on = 'track_id')

# Inspect the resultant dataframe
echo_tracks.head()
```

Out[1]:

	track_id	acousticness	danceability	energy	instrumentalness	liveness	speechiness	tempo	valence	genre_top
0	2	0.416675	0.675894	0.634476	0.010628	0.177647	0.159310	165.922	0.576661	Hip-Hop

1	track_id	acousticness	danceability	energy	instrumentalness	liveness	speechiness	tempo	valence	genre_top
2	5	0.043567	0.745566	0.701470	0.000697	0.373143	0.124595	100.260	0.621661	Hip-Hop
3	134	0.452217	0.513238	0.560410	0.019443	0.096567	0.525519	114.290	0.894072	Hip-Hop
4	153	0.988306	0.255661	0.979774	0.973006	0.121342	0.051740	90.241	0.034018	Rock

## 2. Pairwise relationships between continuous variables

We typically want to avoid using variables that have strong correlations with each other -- hence avoiding feature redundancy -- for a few reasons:

- To keep the model simple and improve interpretability (with many features, we run the risk of overfitting).
- When our datasets are very large, using fewer features can drastically speed up our computation time.

To get a sense of whether there are any strongly correlated features in our data, we will use built-in functions in the `pandas` package.

In [2]:

```
# Create a correlation matrix
corr_metrics = echo_tracks.corr()
corr_metrics.style.background_gradient( )
```

Out[2]:

	track_id	acousticness	danceability	energy	instrumentalness	liveness	speechiness	tempo	valence
track_id	1.000000	-0.372282	0.049454	0.140703	-0.275623	0.048231	-0.026995	-0.025392	0.010070
acousticness	-0.372282	1.000000	-0.028954	-0.281619	0.194780	-0.019991	0.072204	-0.026310	-0.013841
danceability	0.049454	-0.028954	1.000000	-0.242032	-0.255217	-0.106584	0.276206	-0.242089	0.473165
energy	0.140703	-0.281619	-0.242032	1.000000	0.028238	0.113331	-0.109983	0.195227	0.038603
instrumentalness	-0.275623	0.194780	-0.255217	0.028238	1.000000	-0.091022	-0.366762	0.022215	-0.219967
liveness	0.048231	-0.019991	-0.106584	0.113331	-0.091022	1.000000	0.041173	0.002732	-0.045093
speechiness	-0.026995	0.072204	0.276206	-0.109983	-0.366762	0.041173	1.000000	0.008241	0.149894
tempo	-0.025392	-0.026310	-0.242089	0.195227	0.022215	0.002732	0.008241	1.000000	0.052221
valence	0.010070	-0.013841	0.473165	0.038603	-0.219967	-0.045093	0.149894	0.052221	1.000000

## 3. Normalizing the feature data

As mentioned earlier, it can be particularly useful to simplify our models and use as few features as necessary to achieve the best result. Since we didn't find any particular strong correlations between our features, we can instead use a common approach to reduce the number of features called **principal component analysis (PCA)**.

It is possible that the variance between genres can be explained by just a few features in the dataset. PCA rotates the data along the axis of highest variance, thus allowing us to determine the relative contribution of each feature of our data towards the variance between classes.

However, since PCA uses the absolute variance of a feature to rotate the data, a feature with a broader range of values will overpower and bias the algorithm relative to the other features. To avoid this, we must first normalize our data. There are a few methods to do this, but a common way is through *standardization*, such that all features have a mean = 0 and standard deviation = 1 (the resultant is a z-score).

In [3]:

```
# Define our features
features = echo_tracks.drop(['genre_top', 'track_id'], axis = 1 )

# Define our labels
labels = echo_tracks['genre_top']

# Import the StandardScaler
from sklearn.preprocessing import StandardScaler

# Scale the features and set the values to a new variable
```

```
# Scale the features and see the values to a new variable
scaler = StandardScaler()
scaled_train_features = scaler.fit_transform(features)
```

## 4. Principal Component Analysis on our scaled data

Now that we have preprocessed our data, we are ready to use PCA to determine by how much we can reduce the dimensionality of our data. We can use **scree-plots** and **cumulative explained ratio plots** to find the number of components to use in further analyses.

Scree-plots display the number of components against the variance explained by each component, sorted in descending order of variance. Scree-plots help us get a better sense of which components explain a sufficient amount of variance in our data. When using scree plots, an 'elbow' (a steep drop from one data point to the next) in the plot is typically used to decide on an appropriate cutoff.

In [4]:

```
# This is just to make plots appear in the notebook
%matplotlib inline

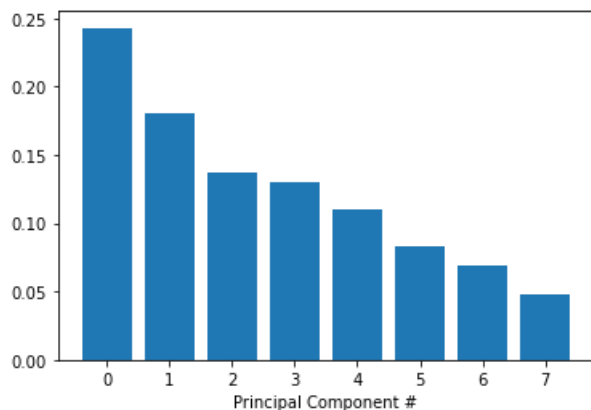
# Import our plotting module, and PCA class
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA

# Get our explained variance ratios from PCA using all features
pca = PCA()
pca.fit(scaled_train_features)
exp_variance = pca.explained_variance_ratio_

# plot the explained variance using a barplot
fig, ax = plt.subplots()
ax.bar(range(8), height = exp_variance)
ax.set_xlabel('Principal Component #')
```

Out[4]:

Text(0.5, 0, 'Principal Component #')



## 5. Further visualization of PCA

Unfortunately, there does not appear to be a clear elbow in this scree plot, which means it is not straightforward to find the number of intrinsic dimensions using this method.

But all is not lost! Instead, we can also look at the **cumulative explained variance plot** to determine how many features are required to explain, say, about 85% of the variance (cutoffs are somewhat arbitrary here, and usually decided upon by 'rules of thumb'). Once we determine the appropriate number of components, we can perform PCA with that many components, ideally reducing the dimensionality of our data.

In [5]:

```
# Import numpy
import numpy as np

# Calculate the cumulative explained variance
```

```

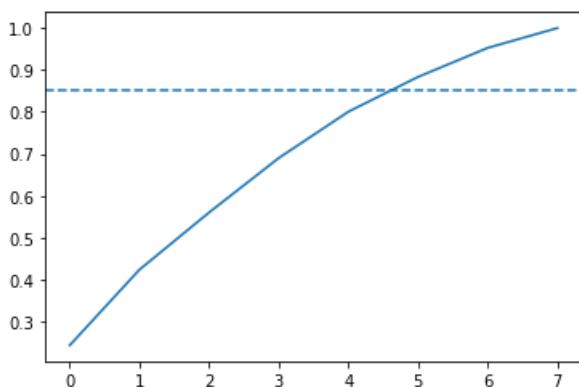
# calculate the cumulative explained variance
cum_exp_variance = np.cumsum(exp_variance)

# Plot the cumulative explained variance and draw a dashed line at 0.85.
fig, ax = plt.subplots()
ax.plot(cum_exp_variance)
ax.axhline(y=0.85, linestyle='--')

# choose the n_components where about 85% of our variance can be explained
n_components = 6

# Perform PCA with the chosen number of components and project data onto components
pca = PCA(n_components, random_state=10)
pca.fit(scaled_train_features)
pca_projection = pca.transform(scaled_train_features)

```

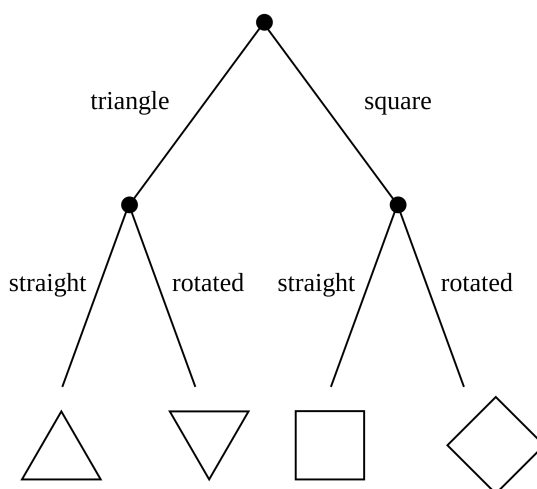


## 6. Train a decision tree to classify genre

Now we can use the lower dimensional PCA projection of the data to classify songs into genres. To do that, we first need to split our dataset into 'train' and 'test' subsets, where the 'train' subset will be used to train our model while the 'test' dataset allows for model performance validation.

Here, we will be using a simple algorithm known as a decision tree. Decision trees are rule-based classifiers that take in features and follow a 'tree structure' of binary decisions to ultimately classify a data point into one of two or more categories. In addition to being easy to both use and interpret, decision trees allow us to visualize the 'logic flowchart' that the model generates from the training data.

Here is an example of a decision tree that demonstrates the process by which an input image (in this case, of a shape) might be classified based on the number of sides it has and whether it is rotated.



In [6]:

```

# Import train_test_split function and Decision tree classifier
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report

# Split our data
train_features, test_features, train_labels, test_labels = train_test_split(pca_projection, labels,

```

```
random_state = 10)
```

```
# Train our decision tree
tree = DecisionTreeClassifier(random_state=10)
tree.fit(train_features, train_labels)

# Predict the labels for the test data
pred_labels_tree = tree.predict(test_features)
classification_report_tree = classification_report(test_labels, pred_labels_tree)
print(classification_report_tree)
```

	precision	recall	f1-score	support
Hip-Hop	0.60	0.60	0.60	235
Rock	0.90	0.90	0.90	966
accuracy			0.84	1201
macro avg	0.75	0.75	0.75	1201
weighted avg	0.84	0.84	0.84	1201

## 7. Compare our decision tree to a logistic regression

Although our tree's performance is decent, it's a bad idea to immediately assume that it's therefore the perfect tool for this job -- there's always the possibility of other models that will perform even better! It's always a worthwhile idea to at least test a few other algorithms and find the one that's best for our data.

Sometimes simplest is best, and so we will start by applying **logistic regression**. Logistic regression makes use of what's called the logistic function to calculate the odds that a given data point belongs to a given class. Once we have both models, we can compare them on a few performance metrics, such as false positive and false negative rate (or how many points are inaccurately classified).

In [7]:

```
# Import LogisticRegression
from sklearn.linear_model import LogisticRegression

# Train our logistic regression and predict labels for the test set
logreg = LogisticRegression(random_state = 10)
logreg.fit(train_features, train_labels)
pred_labels_logit = logreg.predict(test_features)

# Create the classification report for both models
from sklearn.metrics import classification_report as cr
class_rep_tree = cr(test_labels, pred_labels_tree)
class_rep_log = cr(test_labels, pred_labels_logit)

print("Decision Tree: \n", class_rep_tree)
print("Logistic Regression: \n", class_rep_log)
```

Decision Tree:

	precision	recall	f1-score	support
Hip-Hop	0.60	0.60	0.60	235
Rock	0.90	0.90	0.90	966
accuracy			0.84	1201
macro avg	0.75	0.75	0.75	1201
weighted avg	0.84	0.84	0.84	1201

Logistic Regression:

	precision	recall	f1-score	support
Hip-Hop	0.77	0.54	0.64	235
Rock	0.90	0.96	0.93	966
accuracy			0.88	1201
macro avg	0.83	0.75	0.78	1201
weighted avg	0.87	0.88	0.87	1201

## 8. Balance our data for greater performance

## 8. Balance our data for greater performance

Both our models do similarly well, boasting an average precision of 87% each. However, looking at our classification report, we can see that rock songs are fairly well classified, but hip-hop songs are disproportionately misclassified as rock songs.

Why might this be the case? Well, just by looking at the number of data points we have for each class, we see that we have far more data points for the rock classification than for hip-hop, potentially skewing our model's ability to distinguish between classes. This also tells us that most of our model's accuracy is driven by its ability to classify just rock songs, which is less than ideal.

To account for this, we can weight the value of a correct classification in each class inversely to the occurrence of data points for each class. Since a correct classification for "Rock" is not more important than a correct classification for "Hip-Hop" (and vice versa), we only need to account for differences in *sample size* of our data points when weighting our classes here, and not relative importance of each class.

In [8]:

```
# Subset only the hip-hop tracks, and then only the rock tracks
hop_only = echo_tracks.loc[echo_tracks['genre_top'] == 'Hip-Hop']
rock_only = echo_tracks.loc[echo_tracks['genre_top'] == 'Rock']

# sample the rocks songs to be the same number as there are hip-hop songs
rock_only = rock_only.sample(len(hop_only.index), random_state=10)

# concatenate the dataframes rock_only and hop_only
rock_hop_bal = pd.concat([rock_only, hop_only])

# The features, labels, and pca projection are created for the balanced dataframe
features = rock_hop_bal.drop(['genre_top', 'track_id'], axis=1)
labels = rock_hop_bal['genre_top']
pca_projection = pca.fit_transform scaler.fit_transform(features))

# Redefine the train and test set with the pca_projection from the balanced data
train_features, test_features, train_labels, test_labels = train_test_split(pca_projection, labels,
                                                                              random_state=10)
```

## 9. Does balancing our dataset improve model bias?

We've now balanced our dataset, but in doing so, we've removed a lot of data points that might have been crucial to training our models. Let's test to see if balancing our data improves model bias towards the "Rock" classification while retaining overall classification performance.

Note that we have already reduced the size of our dataset and will go forward without applying any dimensionality reduction. In practice, we would consider dimensionality reduction more rigorously when dealing with vastly large datasets and when computation times become prohibitively large.

In [9]:

```
# Train our decision tree on the balanced data
tree = DecisionTreeClassifier(random_state=10)
tree.fit(train_features, train_labels)
pred_labels_tree = tree.predict(test_features)

# Train our logistic regression on the balanced data
logreg = LogisticRegression(random_state=10)
logreg.fit(train_features, train_labels)
pred_labels_logit = logreg.predict(test_features)

# Compare the models
print("Decision Tree: \n", cr(test_labels, pred_labels_tree))
print("Logistic Regression: \n", cr(test_labels, pred_labels_logit))
```

Decision Tree:

	precision	recall	f1-score	support
Hip-Hop	0.74	0.73	0.74	230
Rock	0.73	0.74	0.73	225
accuracy			0.74	455
macro avg	0.74	0.74	0.74	455
weighted avg	0.74	0.74	0.74	455

Logistic Regression:

	precision	recall	f1-score	support
Hip-Hop	0.74	0.73	0.74	230
Rock	0.73	0.74	0.73	225
accuracy			0.74	455
macro avg	0.74	0.74	0.74	455
weighted avg	0.74	0.74	0.74	455

## 10. Using cross-validation to evaluate our models

Success! Balancing our data has removed bias towards the more prevalent class. To get a good sense of how well our models are actually performing, we can apply what's called **cross-validation** (CV). This step allows us to compare models in a more rigorous fashion.

Since the way our data is split into train and test sets can impact model performance, CV attempts to split the data multiple ways and test the model on each of the splits. Although there are many different CV methods, all with their own advantages and disadvantages, we will use what's known as **K-fold** CV here. K-fold first splits the data into K different, equally sized subsets. Then, it iteratively uses each subset as a test set while using the remainder of the data as train sets. Finally, we can then aggregate the results from each fold for a final model performance score.

In [10]:

```
from sklearn.model_selection import KFold, cross_val_score

# Set up our K-fold cross-validation
kf = KFold(n_splits = 10)

tree = DecisionTreeClassifier(random_state=10)
logreg = LogisticRegression(random_state=10)

# Train our models using KFold cv
tree_score = cross_val_score( tree, pca_projection, labels, cv= kf)
logit_score = cross_val_score( logreg, pca_projection, labels, cv = kf)

# Print the mean of each array of scores
print("Decision Tree:", np.mean(tree_score), "Logistic Regression:", np.mean(logit_score))
```

Decision Tree: 0.7489010989010989 Logistic Regression: 0.782967032967033